

Implementación del clasificador Naive Bayes para la Acentuación Automática de Palabras Ambiguas del Español.

Yesenia-N González-Meneses¹, Blanca-Estela Pedroza-Méndez², Francisco López-Briones³, Carlos Pérez-Corona⁴, J-Federico Ramírez-Cruz⁵

INSTITUTO TECNOLÓGICO DE APIZACO.

Av. Instituto Tecnológico s/n Apizaco Tlaxcala México.

¹yeseniaglez@hotmail.com; ²thelismedina@hotmail.com; ³franlopri@hotmail.com;
⁴cperez_corona@hotmail.com; ⁵jframirez@hotmail.com.

Resumen. En este artículo se analiza uno de los problemas más representativos en el tratamiento del lenguaje español que es el de la ambigüedad que existe en la acentuación gráfica de las palabras. En la escritura del lenguaje español existe una marca muy utilizada llamada acento o tilde, esta determina la pronunciación correcta de las palabras. En algunas palabras la tilde siempre se coloca en la misma sílaba, pero hay otras que la pueden llevar o no, o la llevan en diferente sílaba, esto es debido a que estas palabras toman diferentes sentidos dependiendo del contexto donde se encuentran, en este caso se utiliza la llamada tilde diacrítica. La asignación correcta de la tilde diacrítica en este proyecto es vista como un problema de clasificación, donde en base al contexto se determina si las palabras ambiguas llevan esta marca o no. Para lo que se entrenó un modelo con el clasificador Naive Bayes.

Abstract. This paper analyzes one of the most representative problems in treatment of Spanish language, which is the ambiguity that exists in graphic accentuation of words. In writing of the Spanish language is widely used a mark called accent or tilde, this determines the correct pronunciation of words. In some words the accent is always placed in the same syllable, but there are others that can lead or not, or can lead in a different syllable, this is because these words take on different meanings depending on the context where they are, in this case is used the diacritical tilde. In this project the correct allocation of diacritical tilde is seen as a classification problem, where the context determines whether ambiguous words lead the mark or not. For this we trained and tested a model with the Naive Bayes classifier.

Keywords: Ambigüedad en la acentuación; Clasificador Naive Baye, Etiquetado de texto.

1 Introducción.

El Procesamiento del Lenguaje Natural o PLN es un área de la Inteligencia Artificial, dependiente directamente de la Lingüística Computacional. Así mismo, el PLN es un

componente importante de las interfaces de usuarios y los sistemas inteligentes y uno de los objetivos que persigue es el perfecto análisis y entendimiento de los lenguajes humanos [18].

Los esfuerzos de investigación en PLN han sido dirigidos hacia tareas intermedias que dan sentido a alguna de las múltiples características estructurales inherentes a los lenguajes, sin requerir un entendimiento completo. Una de esas tareas es la asignación de categorías gramaticales o morfosintácticas (sustantivo, adjetivo, verbo, etc.) a cada una de las palabras de una oración. Este proceso se denomina también etiquetación [26]. El proceso de etiquetación debe eliminar ambigüedades y encontrar cual es el papel más probable que juega cada palabra dentro de una frase. Dicho proceso debe ser capaz también de asignar una etiqueta a cada una de las palabras que aparecen en un texto, y garantizar de alguna manera que esa es la etiqueta correcta.

El problema más difícil que se enfrenta en el procesamiento del lenguaje es la ambigüedad: que es cuando pueden admitirse distintas interpretaciones a partir de la representación o cuando existe confusión al tener diversas estructuras y no tener los elementos necesarios para eliminar las incorrectas [17]. Este problema se presenta en todos los niveles del lenguaje, sin excepción [13]. Desde el nivel morfológico (palabras), hasta el discurso (o pragmática).

1.1 Descripción del problema.

En el proceso de escribir textos en lenguaje español, muchas veces cometemos errores ortográficos, debido a que es muy común que la gente olvide como utilizar las reglas del idioma que regulan esta tarea. Aunque estas reglas son enseñadas desde niños, se van olvidando ya que no se pone el suficiente empeño en aplicarlas, una de las causas de esto es que como humanos tenemos la inteligencia suficiente para entender los textos sin importar que estos no estén escritos correctamente. Otro de los problemas es que el idioma español en sí es muy ambiguo en la escritura de las palabras, por ejemplo hay palabras que tienen idéntica pronunciación pero su escritura y su significado es diferente (*palabras homófonas*), ejemplo: *tuvo* y *tubo*, *huno* y *uno*. Otro caso es el de la *polisemia*, que es cuando una palabra tiene diferentes significados, por ejemplo la palabra *banco* que puede tener significado de institución de crédito, o de asiento sin respaldo, etc. En este caso lo que permite darle el sentido correcto a la palabra, es el contexto en el que se encuentra.

Otra de las cosas que genera ambigüedad en el idioma español, es la acentuación gráfica de las palabras, ya que existen palabras que se escriben igual pero pueden o no llevar acento dependiendo del contexto de la frase que contiene la palabra. Por ejemplo a la palabra “gráfica” se le debe colocar acento en la “a” de la sílaba “gra” si la palabra es un sustantivo, pero si la palabra dentro de la frase se maneja como un verbo, la sílaba tónica es “fi” y de acuerdo a las reglas de acentuación no lleva acento, ya que es una palabra grave que termina en vocal. Por tanto podemos observar que existe una relación entre la acentuación y las etiquetas morfosintácticas que se le asignan a las palabras.

En este artículo nos enfocamos al análisis de las reglas de acentuación y se propone un modelo basado en métodos de aprendizaje automático, aplicando el clasificador Naive Bayes para dar solución al problema de la ambigüedad al asignar el acento diacrítico. El clasificador analiza el contexto de la frase en base a las etiquetas morfosintácticas asignadas a las palabras y determina cuando una palabra debe o no llevar acento diacrítico, para lo cual se deben primero corregir las palabras con acento gráfico, esto es para disminuir el número de errores por omisión de acentos y al mismo tiempo para que las etiquetas generadas sean más precisas. El diccionario utilizado en este proyecto se generó como una de las etapas iniciales, donde se identificaron además otro tipo de ambigüedades en la acentuación gráfica de palabras; para la fase de la etiquetación se utilizó el módulo para este fin del paquete Freeling [19].

La Desambiguación del Sentido de las Palabras WSD (Word Sense Disambiguation) es en esencia una tarea de clasificación: Los sentidos de las palabras son las clases, el contexto provee la evidencia y cada una de las palabras es asignada a una o más de las posibles clases basado en la evidencia [12]. El clasificador Naive Bayes es uno de los algoritmos que estiman probabilidades a posteriori. Este clasificador asume, para una muestra x , que sus atributos x_1, x_2, \dots, x_n presentan una independencia condicional dado el valor de la clase, por lo que la probabilidad condicional puede expresarse como el producto de funciones de probabilidad condicional de cada atributo por separado [22]. En este sentido los atributos utilizados para la desambiguación son las palabras en contexto a la palabra ambigua, siendo los valores de cada atributo las etiquetas morfosintácticas asignadas por Freeling, para que, calculando probabilidades por cada uno de estos valores con respecto a la clase de salida, se pueda definir la clase a la que pertenecen dichas palabras.

2 El lenguaje Español.

La ortografía es la rama de la gramática que se ocupa de la escritura correcta [14]. Según el diccionario de la Real Academia Española se define como: “Conjunto de normas que regulan la escritura de una lengua”.

Dentro del lenguaje español todas las palabras tienen una sílaba que se pronuncia con mayor intensidad, esto es lo que se conoce como acento prosódico, que es el mayor relieve con que se pronuncia una determinada sílaba dentro de una palabra. Otro tipo de acento que se maneja dentro del español es el acento gráfico u ortográfico, que es el signo con el cual, en determinados casos, se representa en la escritura el acento prosódico [16].

En las reglas de la gramática del español se hace una clasificación para los acentos como sigue [15]:

- **Tilde diacrítica o acento diacrítico.**- Es la marca que se coloca sobre alguna de las vocales dentro de una palabra para permitir diferenciar entre los significados de ésta.

- **Acento gráfico.**- Esta no se utiliza para diferenciar entre los significados sino para saber la pronunciación correcta de una palabra, en el caso contrario la colocación de esta marca la define la pronunciación de la palabra.

El error más común cuando escribimos textos, es la omisión tanto del acento gráfico como del acento diacrítico, ya que aunque no es difícil identificar la sílaba tónica, sí lo es recordar las reglas. Actualmente, es muy común el uso de procesadores de texto, que ya tienen incluido un diccionario de palabras para ayudar a la acentuación, pero cuando se trata de palabras con ambigüedad en la acentuación, el procesador no indica si deben o no llevar acento.

3 Clasificación.

La clasificación es el punto principal en esta investigación, ya que la asignación de la tilde diacrítica a las palabras ambiguas se modela como un problema de clasificación, donde las clases para cada palabra es si lleva o no lleva la tilde.

La clasificación es la tarea de aproximar una función objetivo desconocida $\Phi : I \times C \rightarrow \{T, F\}$ por medio de una función $\Theta : I \times C \rightarrow \{T, F\}$ llamada clasificador, donde $C = \{c_1, c_2, \dots, c_{|C|}\}$ es un conjunto de clases definido, e I es un conjunto de instancias del problema. Cada instancia $ij \in I$ es representada como una lista $A = \{a_1, a_2, \dots, a_{|A|}\}$ de valores característicos, conocidos como atributos. i.e. $ij = \{a_{1j}, a_{2j}, \dots, a_{|A|j}\}$. Si $\Phi : I \times C \rightarrow T$ entonces ij es llamado un ejemplo positivo de c_i , mientras que si $\Theta : I \times C \rightarrow F$ es llamado un ejemplo negativo de c_i [23]. En general no se conoce la descripción exacta de las muestras, por lo que el sistema es entrenado a priori para ajustarse a las características propias del problema. A este proceso de adquirir e integrar conocimiento a un sistema de clasificación a partir de ejemplos, se le conoce como aprendizaje o entrenamiento [22].

3.1 Clasificador Naive Bayes

Uno de los métodos supervisados que estiman probabilidades a posteriori es el algoritmo Naive Bayes. Este clasificador asume, para una muestra x , que sus atributos x_1, x_2, \dots, x_n presentan una independencia condicional dado el valor de la clase, por lo que la probabilidad condicional puede expresarse como el producto de funciones de probabilidad condicional de cada atributo por separado.

(1)

$$p(x|\omega_i) = \prod_{j=1}^n p(x_j|\omega_i) \quad (1)$$

Usando el teorema de Bayes, la probabilidad a posteriori se escribe como,

$$p(\omega_i|x) = p(\omega_i) \prod_{j=1}^n p(x_j|\omega_i) \quad (2)$$

Finalmente, el algoritmo Naive Bayes asigna una muestra \mathbf{x} a una de las L clases existentes utilizando la función:

$$\omega^* = \operatorname{argmax}_{\omega_j} p(\omega_j) \prod_{i=1}^n p(x_i | \omega_j) \quad (3)$$

3.2 Validación cruzada (Cross validation).

La validación cruzada. Conocido como método π o rotación, genera aleatoriamente una partición en K bloques de tamaño N/K . El entrenamiento (training) se lleva a cabo empleando $K - 1$ bloques, mientras que el subconjunto restante es empleado como prueba (test). Este procedimiento es repetido K veces, eligiendo en cada iteración una parte diferente para prueba. Una extensión a este método es el llamado stratified cross validation (validación cruzada estratificada) con el que, para cada partición, las clases se encuentran distribuidas según sus probabilidades a priori en el conjunto original. Por otra parte, para una mejor estimación, el proceso es repetido P veces. La Fig. 1 muestra un ejemplo de validación cruzada con $K = 3$ [24].

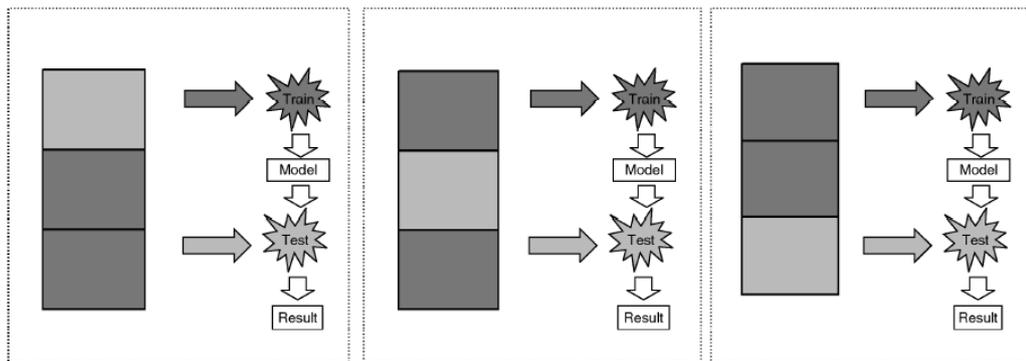


Fig. 1. Validación cruzada, $K = 3$.

Para este trabajo se utilizó la validación cruzada estratificada con $K = 10$.

3.3 Evaluación de la efectividad del clasificador.

Las métricas de evaluación más empleadas para medir la efectividad de los clasificadores son la tasa de errores y la tasa de aciertos. Estas, para un problema de dos clases, pueden obtenerse a partir de una matriz de confusión Tabla. 1.

Tabla 1. Matriz de confusión para un problema de dos clases

	Positivos (total)	Negativos (total)
Positivos (clasificador)	Verdaderos Positivos (VP)	Falsos Positivos (FP)
Negativos (clasificador)	Falsos Negativos (FN)	Verdaderos Negativos(VN)

Estas tasas pueden calcularse como:

$$Acc = \frac{vp + vn}{vp + fn + vn + fp} \quad (4)$$

y

$$Err = \frac{fp + fn}{vp + fn + vn + fp} \quad (5)$$

Aunque estas medidas no resultan apropiadas debido a que no consideran distintos tipos de errores, ya que se muestran fuertemente sesgadas a favor de la clase mayoritaria. Por ejemplo, considerando un problema binario cuya clase positiva contiene un 1% de objetos sobre el conjunto total; en tal situación, una simple estrategia de asignar todas las muestras a la clase negativa ofrecería una tasa de aciertos del 99%, sin embargo, tal clasificador carecería de valor alguno [22]. Lo que ha motivado a la búsqueda de medidas alternativas. Algunos ejemplos son los siguientes:

- Tasa de Verdaderos Positivos (Sensibilidad), es el porcentaje de ejemplos positivos que son correctamente clasificados.

$$tvp = \frac{vp}{vp + fn} \quad (6)$$

- Tasa de Verdaderos Negativos (Especificidad), es el porcentaje de ejemplos negativos que son clasificados como positivos.

$$tvn = \frac{vn}{cn + fp} \quad (7)$$

- Tasa de Falsos Positivos, es el porcentaje de ejemplos negativos que son erróneamente clasificados.

$$tfp = \frac{fp}{vn + fp} \quad (8)$$

- Tasa de Falsos Negativos, es el porcentaje de ejemplos positivos que son clasificados como negativos.

$$tfn = \frac{fn}{vp + fn} \quad (9)$$

- Precisión, se define como el porcentaje de ejemplos que fueron etiquetados correctamente como positivos, con respecto a todas las muestras que fueron etiquetadas como tal.

$$precision = \frac{vp}{vp + fp} \quad (10)$$

Curvas ROC

La curva ROC es una metodología de análisis desarrollada por ingenieros eléctricos y de radar durante la Segunda Guerra Mundial, con la finalidad de resolver problemas prácticos en la detección de señales. El espacio de la curva ROC es un gráfico bidimensional que permite visualizar, organizar y seleccionar clasificadores basados en su efectividad, en nuestro caso se utilizará para comparar los diferentes parámetros para determinar con cuales se obtienen mejores resultados. Mediante esta representación es posible conocer la relación entre los Verdaderos Positivos y los Falsos Negativos. La Fig. 2 muestra un espacio ROC, cuyo eje Y representa la sensibilidad y el eje X la especificidad. En esta misma figura, se encuentran cinco clasificadores etiquetados de la A a la E [22].

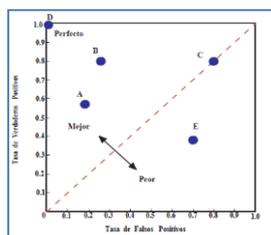


Fig. 2. Espacio de una curva ROC.

4. Metodología y desarrollo

En esta sección se describe la metodología utilizada en la realización del proyecto, explicando cada una de las etapas y los pasos realizados en ellas. Los programas desarrollados en cada una de estas etapas se realizaron en Java, utilizando el entorno de desarrollo Integrado NetBeans IDE 6.8, y Excel (Macros). En la Fig. 3 se muestra el diagrama general del proyecto, del cual se destacan las actividades representadas en los módulos de la Fig. 4 enumerados de acuerdo al orden en que se fueron realizando.

Las actividades de la figura 4 son explicadas a continuación:

- 1) **Clasificación de palabras con acento (tilde) de acuerdo a las reglas que determina la Real Academia Española (RAE).** Al realizar un análisis detallado de las reglas se identificaron los tipos de palabras que contienen tilde por lo que se presenta una clasificación y las diferentes formas en que se pueden solucionar.
- 2) **Diccionario de palabras con acento gráfico.** De acuerdo a las reglas expuestas en el capítulo anterior se pueden ver principalmente dos clasificaciones de palabras con acento, las que tienen *acento gráfico* y las que tienen *acento diacrítico*. Por las definiciones dadas a cada uno de estos acentos se puede ver que el *acento diacrítico* se utiliza para diferenciar entre significados de las palabras, mientras que el acento gráfico no presenta ambigüedad en su significado. Por lo que la generación de un diccionario con su forma correcta es suficiente para la corrección de este tipo de palabras.
- 3) **Palabras con acento diacrítico a tratar con Naive Bayes.** El principal problema que se abordó en este proyecto es el de la asignación correcta del acento diacrítico, por lo que el primer paso fue identificar las palabras que lo Necesitan y determinar la forma en que se llevó a cabo la clasificación de acuerdo al análisis de frases con palabras ambiguas.
- 4) **Obtención de ejemplos.** En este módulo se obtuvieron ejemplos para cada una de las formas que puede tomar cada palabra ambigua, los ejemplos se extrajeron del banco de datos CREA, disponible en línea en <http://corpus.rae.es/creanet.html>.
- 5) **Pre-procesamiento de ejemplos.** Partiendo del planteamiento del problema, donde se dice que la omisión de acentos es uno de los principales errores en la escritura y el problema a corregir en este proyecto, se eliminan todos los

acentos contenidos en los ejemplos, para posteriormente colocarlos a las palabras que les corresponda.

- 6) **Corrección de palabras con acento gráfico.** El diccionario obtenido del módulo dos se aplicará en esta parte, que es la de restauración de acentos a las palabras de esta clase.

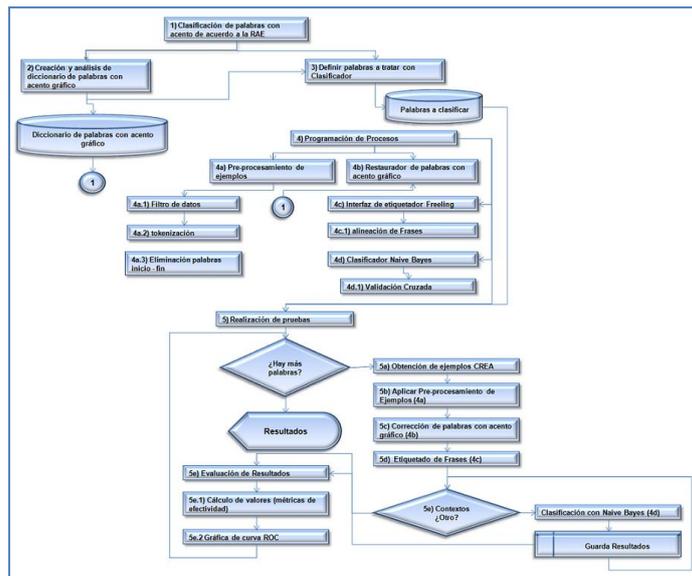


Fig. 3 Diagrama General del Proyecto.

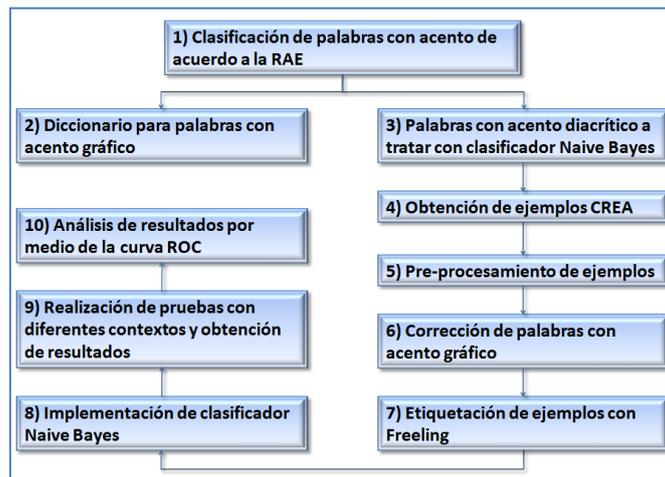


Fig. 4. Actividades realizadas durante el proyecto

- 7) **Etiquetación de ejemplos con Freeling.** Las posibles combinaciones de palabras para formar frases dentro del lenguaje es un número infinito dado que la cantidad de palabras es muy grande, sin embargo siguen una misma estructura definida por la gramática del idioma en base a categorías gramaticales (etiquetas), por este motivo se realiza una etiquetación para obtener características de las palabras y clasificar en base a esa información.
- 8) **Implementación del clasificador Naive Bayes.** Este módulo es el más importante de todo el proyecto, es donde se le asigna el sentido correcto a la palabra ambigua en base a la información contenida en las etiquetas que regresa Freeling, los resultados que regrese el clasificador son evaluados por medio de la validación cruzada, que evalúa principalmente la capacidad de generalización del modelo entrenado.
- 9) **Realización de pruebas con diferentes contextos y obtención de resultados.** Se realizaron diferentes pruebas tomando en cuenta contextos variados, tomando como máxima referencia tres etiquetas hacia adelante de la palabra ambigua, tres etiquetas hacia atrás y la etiqueta de la palabra ambigua.
- 10) **Análisis de resultados por medio de la curva ROC.** Los resultados regresados por el clasificador pueden ser vistos como una matriz de confusión, de la cual se pueden obtener los valores necesarios para analizarlos por medio de este método, y así determinar el mejor contexto asociado a cada palabra para desambiguarla.

Pruebas y Resultados

Para la mayoría de palabras se realizaron cuatro clases de pruebas: tomando en cuenta que las clases están equilibradas, es decir, clasificando de tal manera que se tenga una probabilidad del 50% (Tabla. 2) ser de una clase o de otra, esto para ver comportamiento del clasificador; y tomando en cuenta la proporción de acuerdo a las consultas realizadas en el corpus CREA (Tabla. 3). Estas dos formas a su vez fueron divididas en dos, tomando en cuenta la palabra ambigua y sin tomarla en cuenta, dado que en estos ejemplos la palabra ambigua es etiquetada de diferentes maneras dependiendo del contexto, pero inclinándose hacia una de las dos clases. En los casos como (*mi, te, tu, cuan* y *quien*) solo se realizaron pruebas sin tomar en cuenta la palabra ambigua, ya que toma la misma proporción que la clase.

En las tablas 2 y 3 se presenta un ejemplo de la forma en que se fueron realizando las pruebas, donde los valores mostrados son explicados a continuación:

- **Proporción:** Distribución de los datos en pruebas con el clasificador.
- **Contexto:** Las palabras circundantes a la palabra ambigua (desde -3 *amb* $+3$; hasta -3 $+3$).
- **Acc** (exactitud): Porcentaje de ejemplos clasificados correctamente, definido por la ecuación (2.4)
- **VN** (Verdaderos Negativos): Ejemplos clasificados correctamente como ejemplos sin acento.
- **FP** (Falsos Positivos): Ejemplos clasificados incorrectamente como ejemplos sin acento.

- **VP** (Verdaderos Positivos): Ejemplos clasificados correctamente como ejemplos con acento)
- **FN** (Falsos Negativos): Ejemplos clasificados incorrectamente como ejemplos con acento.

Tabla 2. Pruebas proporción 50 - 50

contexto	Proporcion 50 - 50				
	Acc	VN	FP	VP	FN
amb +3	0,7557	33,9	10,1	32,6	11,4
-1 amb +3	0,8966	38,3	5,7	40,6	3,4
-2 amb +3	0,9216	40,3	3,7	40,8	3,2
-3 amb +3	0,9045	39,2	4,8	40,4	3,6
amb +2	0,7557	34,3	9,7	32,2	11,8
-1 amb +2	0,9148	40,4	3,6	40,1	3,9
-2 amb +2	0,9136	40,0	4,0	40,6	3,6
amb +1	0,6727	31,0	13,0	28,2	15,8
-2 amb +1	0,9136	40,0	4,0	40,4	3,6
-3 amb +1	0,8909	37,8	6,2	40,6	3,4
-1 amb	0,8977	39,4	4,6	39,6	4,4
-3 amb	0,8784	37,6	6,4	39,7	4,3
+3	0,7557	33,6	10,4	32,9	11,1
-1 +3	0,8886	38,7	5,3	39,5	4,5
-2 +3	0,9057	39,6	4,4	40,1	3,9
+2	0,7511	34,2	9,8	31,9	12,1
-2 +2	0,9170	40,6	3,4	40,1	3,9
-3 +2	0,8955	38,8	5,2	40,0	4,0
+1	0,6773	29,2	14,8	30,4	13,6
-1 +1	0,8989	39,4	4,6	39,7	4,3
-3 +1	0,9045	39,4	4,6	40,2	3,8
-1	0,8795	39,2	4,8	38,2	5,8
-2	0,8909	39,1	4,9	39,3	4,7

En estos ejemplos están marcados los mejores resultados de acuerdo a las proporciones que se tomaron en cuenta, siendo el valor de referencia la exactitud. En la Tabla 2 para la proporción 50 – 50 la exactitud llega al 92.16% (contexto -2 amb +3), mientras que en la Tabla .2 (proporción 97 - 03) la exactitud supera el valor mayor de la proporción (97%) con un valor del 98.48% (contexto -2 +2).

Los siguientes son los valores que aparecerán como columnas, además de las anteriores, en las tablas de las pruebas por cada una de las palabras. Estos valores son las métricas utilizadas para el análisis de resultados:

TVP (Tasa De Verdaderos Positivos): Porcentaje de ejemplos positivos que son correctamente clasificados, definido por la ecuación (6)

- **TVN** (Tasa de Verdaderos Negativos): Porcentaje de ejemplos negativos que son clasificados como positivos, definido por la ecuación (7)
- **TFP** (Tasa de Falsos Positivos): Porcentaje de ejemplos negativos que son erróneamente clasificados, definido por la ecuación (8)
- **AUC** (Área Bajo La Curva): área bajo la curva ROC.
- Donde los valores que son graficados en la curva ROC son el **TVP** y el **TFP**, siendo el área marcada el valor de **AUC**.

- Dentro de la investigación se hizo un análisis con todos estos parámetros para diferentes clases de palabras, a continuación se muestra un ejemplo para palabras con terminación *-o* (sustantivo / verbo).
- La Tabla 4 presenta los mejores resultados para las palabras con terminación *-o* en las diferentes pruebas. En la Fig. 4 se pueden ver gráficamente estos resultados.
- Esta prueba es la más importante dentro del proyecto, ya que se está demostrando que el trabajar con etiquetas no solo permite generalizar las palabras en contexto a la palabra ambigua como en las pruebas anteriores, sino que también es posible utilizar etiquetas para generalizar palabras ambiguas, en este caso los verbos que, como se mencionó en el capítulo anterior, se están probando diez palabras diferentes como si fueran una sola, esto al ser de las mismas características.

Table 3. Pruebas proporción CREA

contexto	Proporción 97 - 03				
	Acc	VN	FP	VP	FN
amb +3	0,9630	44,2	0,8	0,1	0,9
-1 amb +3	0,9739	44,4	0,6	0,4	0,6
-2 amb +3	0,9783	44,7	0,3	0,3	0,7
-3 amb +3	0,9761	44,7	0,3	0,2	0,8
amb +2	0,9761	44,6	0,4	0,3	0,7
-1 amb +2	0,9826	44,8	0,2	0,4	0,6
-2 amb +2	0,9804	44,7	0,3	0,4	0,6
amb +1	0,9717	44,7	0,3	0,0	1,0
-2 amb +1	0,9826	44,9	0,1	0,3	0,7
-3 amb +1	0,9783	44,7	0,3	0,3	0,7
-1 amb	0,9804	44,8	0,2	0,3	0,7
-3 amb	0,9717	44,7	0,3	0,0	1,0
+3	0,9587	44,1	0,9	0,0	1,0
-1 +3	0,9826	44,8	0,2	0,4	0,6
-2 +3	0,9804	44,7	0,3	0,4	0,6
+2	0,9674	44,5	0,5	0,0	1,0
-2 +2	0,9848	44,9	0,1	0,4	0,6
-3 +2	0,9804	44,7	0,1	0,4	0,6
+1	0,9783	45,0	0,0	0,0	1,0
-1 +1	0,9783	44,6	0,4	0,4	0,6
-3 +1	0,9783	44,8	0,2	0,2	0,8
-1	0,9826	44,9	0,1	0,3	0,7
-2	0,9804	44,8	0,2	0,3	0,7

Tabla 4. Resultados para sustantivo/verbo “palabras con terminación o”

	proporción	contexto	Acc	VN	FP	VP	FN	tvp	tvn	tfp	AUC
1	50 - 50	-1 amb	0,9444	46,8	2,7	46,7	2,8	0,9434	0,9455	0,0545	0,9444
2	50 - 50	-2	0,9384	46,1	3,4	46,8	2,7	0,9455	0,9313	0,0687	0,9384
3	91 - 09	-2 amb +1	0,9680	45,0	0,5	3,4	1,1	0,7556	0,9890	0,0110	0,8723
4	91 - 09	-2 +1	0,9680	45,1	0,4	3,3	1,2	0,7333	0,9912	0,0088	0,8623

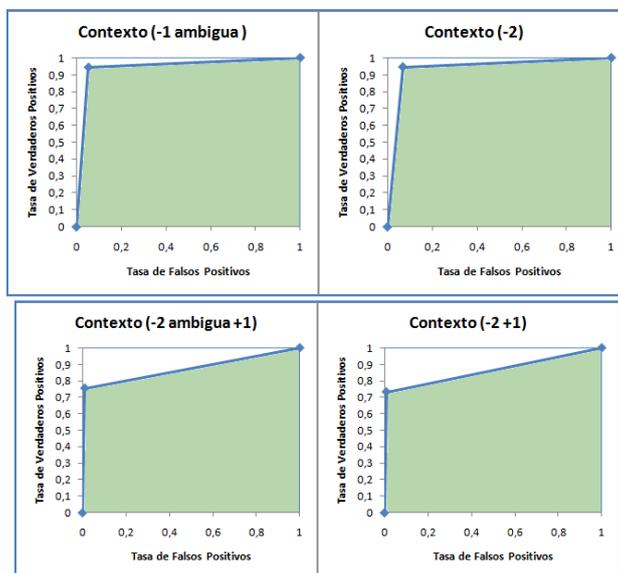


Fig. 4 Área bajo la curva ROC para palabras terminación -o

Conclusiones

Los resultados obtenidos de las diferentes palabras con acento diacrítico fueron buenos, con una exactitud que va desde el 72.12% (demostrativo *cuan*) hasta un 98.94% (monosílabo *se*) cuando se toman en cuenta clases balanceadas. Y tomando en cuenta clases desbalanceadas (proporción CREA) un valor AUC (área bajo la curva ROC) que va desde un 67.77% (demostrativo *cuan*) hasta un 96.15% (monosílabo *te*).

Los resultados más bajos que se obtuvieron fueron para el interrogativo *cuan*, los cuales se dieron debido a que en el corpus CREA, de donde se obtuvieron los datos para el proyecto, no contenía muchos ejemplos para esta palabra, lo que nos dice que no es muy común su uso y por lo mismo en algunos de los ejemplos están acentuadas incorrectamente.

Otra de las cosas que se puede concluir, es que para los monosílabos un contexto cercano es suficiente para desambiguar, mientras que para los interrogativos es necesario un contexto mayor.

Referencias

- [1] Martínez, D. (2008). "On the use of Automatically Acquired Examples for All-Nouns Word Sense Disambiguation". *Journal of Artificial Intelligence Research* 33.
- [2] Tejada, J. (2006). "Desambiguación de Sentidos de Palabras usando relaciones sintácticas como contexto local". *MICAI*

- [3] Suárez, A. (2001). “*Estudio de cooperación de métodos de desambiguación léxica: Marcas de Especificidad vs. Máxima Entropía*”. Procesamiento del lenguaje natural. N° 27
- [4] Universidad Carlos III (consultado 2008) “*Procesamiento del Lenguaje Natural para Recuperación de Información*” <http://pln-ri-hmm.orgfree.com/pln.html>
- [5] Pancardo, A. (2004) “*Desambiguación Léxica de Sustantivos usando la Web*”. Workshop on Lexical Resources and the Web for Word Sense Disambiguation. IX Ibero-American Conference on Artificial Intelligence IBERAMIA
- [6] Carbonell, J. (1992) “*El procesamiento del lenguaje natural, tecnología en transición*”. Carnegie Mellon University.
- [7] Fernández, S. (2006) “*Nueva Propuesta de Desambiguación de Sentidos de Palabras para nombres en un sistema de Búsqueda de Respuestas*”. Procesamiento del lenguaje natural. N° 36
- [8] Jordi Atserias i Batalla. (2006) “*Un Enfoque Integrado para la Desambiguación*”. Procesamiento del lenguaje natural, N° 35
- [9] Villegas, M. (1998) “*El léxico PAROLE del Español*”. Procesamiento del lenguaje natural. N° 23
- [10] Crandall, D. (2005) “*Automatic accent restoration in Spanish text*”. Spring 2005 course project for CS 674.
- [11] Montiel, R. (2010) “*Propuesta de un modelo para la acentuación automática de palabras ambiguas del español, utilizando etiquetado de texto*” Programación Matemática y Software. Vol.2. Num.1
- [12] Ríos Gaona, M. (2008). “*Desambiguación de sentidos de palabras usando sinónimos*”. ESCOM-IPN.
- [13] Traductores. Capítulo 1. Lenguajes. (Consultado Junio 2009) <http://tikal.cifn.unam.mx/~jsegura/academic/traductores/Cap1.htm>.
- [14] Miguel Ángel Monjas Llorente. (Consultado junio 2009) “*Cómo acentuar en español*”. Versión 2.01. 2 de febrero de 1998 <http://www.dat.etsit.upm.es/~mmonjas/acentos.html>
- [15] Real Academia Española (1999) “*ORTOGRAFÍA de la LENGUA ESPAÑOLA*”. Edición revisada por las Academias de la lengua Española.
- [16] Real Academia Española (2005) “*DICCIONARIO PANHISPÁNICO DE DUDAS*” Primera Edición.
- [17] Gelbukh. A. Galicia Haro. S. (2007) “*INVESTIGACIONES EN ANÁLISIS SINTÁCTICO PARA EL ESPAÑOL*”. Instituto Politécnico Nacional. Primera edición.
- [18] Moreno Sandoval. A. (1998) “*LINGÜÍSTICA COMPUTACIONAL: Introducción a los modelos simbólicos, estadísticos y biológicos*”. MADRID, ESPAÑA, SINTESIS.
- [19] Universitat Politècnica de Catalunya (consultado noviembre de 2010) “*Freeling Home Page*”. <http://nlp.lsi.upc.edu/freeling/> Centro de investigación TALP, Universitat Politècnica de Catalunya.
- [20] Mitchel, T.(1997) “*Machine Learning*” McGaw Hill.
- [21] Christopher D. Manning and Hinrich S. (1999) “*Foundations of statistical Natural Language Processing*” Second Printing. The MIT Press Cambridge, Massachusetts.
- [22] Garcia, V. (2010) “*Distribuciones de Clases No Balanceadas: Métricas, Análisis de Complejidad y Algoritmos de Aprendizaje*”. Tesis Doctoral. Departament de llenguatges i Sistemes Informàtics, Universitat Jaume I.
- [23] Sánchez, C. R. (2008) “*Clasificación de Entidades Nombradas utilizando Información Global*” Tesis de Maestría, INAOE.
- [24] Refaeilzadeh, Payam. (2008) “*Cross-Validation*”. Arizona State University.

- [25] Montiel, R. (2009) "*Etiquetación de frases en español para la acentuación de palabras con acento diacrítico mediante método híbrido, considerando el contexto en cuestión*" Tesis de Maestría. Instituto Tecnológico de Apizaco.
- [26] Simard M. (1996). "*Automatic Restoration of Accents In French Text*". Industry Canada. Centre for Information Technology Innovation (CITI). *Automatic Restoration*